

# A rare event approach to high dimensional Approximate Bayesian computation

Dennis Prangle, Richard G. Everitt and Theodore Kypraios

**Abstract** Approximate Bayesian computation (ABC) methods permit approximate inference for intractable likelihoods when it is possible to simulate from the model. However they perform poorly for high dimensional data, and in practice must usually be used in conjunction with dimension reduction methods, resulting in a loss of accuracy which is hard to quantify or control. We propose a new ABC method for high dimensional data based on rare event methods which we refer to as RE-ABC. This uses a latent variable representation of the model. For a given parameter value, we estimate the probability of the rare event that the latent variables correspond to data roughly consistent with the observations. This is performed using sequential Monte Carlo and slice sampling to systematically search the space of latent variables. In contrast standard ABC can be viewed as using a more naive Monte Carlo estimate. We use our rare event probability estimator as a likelihood estimate within the pseudo-marginal Metropolis-Hastings algorithm for parameter inference.

We provide asymptotics showing that RE-ABC has a lower computational cost for high dimensional data than standard ABC methods. We also illustrate our approach empirically, on a Gaussian distribution and an application in infectious disease modelling.

**Keywords** ABC, Markov chain Monte Carlo, sequential Monte Carlo, slice sampling, infectious disease modelling

# 1 Introduction

Approximate Bayesian computation (ABC) is a family of methods for approximate inference, used when likelihoods are impossible or impractical to evaluate numerically but simulating datasets from the model of interest is straightforward. ABC can be viewed as a *nearest neighbours* method. It simulates datasets given various parameter values, and finds the closest matches, in some sense, to the observed dataset. The corresponding parameters are used as the basis for inference. Various Monte Carlo methods have been adapted to implement this idea, including rejection sampling (Beaumont et al., 2002), Markov chain Monte Carlo (MCMC) (Marjoram et al., 2003) and sequential Monte Carlo (SMC) (Sisson et al., 2009). However it is well known that nearest neighbours approaches becomes less effective for higher dimensional data, a phenomenon referred to as the *curse of dimensionality*. The problem is that even under the best parameter values, it is rare for a high dimensional simulation to match a fixed target well, essentially because there are many random components all of which must be close matches to observations.

In this paper we propose a method to deal with this issue and permit high dimensional ABC. The idea involves introducing latent variables  $x$ . We assume data is a deterministic function  $y(\theta, x)$ , where  $\theta$  is a vector of parameters. Hence  $x$  encapsulates all the randomness which occurs in the simulation process. Our approach is, for a particular  $\theta$  value, to use rare event methods to estimate the probability of  $x$  values occurring which produce  $y(\theta, x) \approx y_{\text{obs}}$ . As discussed later, this probability equals, up to proportionality, the approximate likelihood of  $\theta$  used in existing ABC algorithms. We estimate this probability using SMC algorithms for rare events from Cérou et al. (2012). The resulting estimates are unbiased or low bias, depending on the algorithm, and can be used by many inference methods. We concentrate on the pseudo-marginal Metropolis Hastings algorithm (Andrieu and Roberts, 2009), which outputs a sample from a distribution approximating the Bayesian posterior.

The intuition for the rare event probability estimates we use is as follows. Given  $\theta$ , standard ABC methods effectively simulate one or several  $x$  values from their prior and calculate

a Monte Carlo estimate of  $\Pr(y(\theta, x) \approx y_{\text{obs}})$ . This relative error of this estimate has high variance when the probability is small, as is the case when we require close matches. The rare event technique of *splitting* uses nested sets of latent variables  $A_1 \supset A_2 \supset \dots \supset A_T$ , representing increasingly close matches. We aim to estimate  $\Pr(A_1)$ ,  $\Pr(A_2|A_1)$ ,  $\Pr(A_3|A_2), \dots$  and take the product. If these probabilities are all relatively large then the variance of the final estimator’s relative error is smaller than using a single stage of Monte Carlo (For a crude variance analysis justifying this see L’Ecuyer et al., 2007. Cérou et al., 2012 prove more detailed results for their SMC algorithms which we summarise later.) We can estimate  $\Pr(A_1)$  using Monte Carlo with  $N$  samples. Next we reuse the  $x$  samples with  $x \in A_1$ . We sample randomly from these  $N$  times and, to avoid duplicates, perturb each appropriately. We found a good perturbation method was a slice sampling algorithm from Murray and Graham (2016). The resulting sample is used to find a Monte Carlo estimate of  $\Pr(A_2|A_1)$ . We carry on similarly to estimate the remaining conditional probabilities.

## 1.1 Related literature

The most popular approach to deal with the curse of dimensionality in ABC is *dimension reduction*. Here high-dimensional datasets are mapped to lower dimensional vectors of features, often referred to as *summary statistics*. The quality of a match between simulated and observed data is then judged based only on their corresponding summary vectors. However, using summary statistics involves some loss of information about the posterior which is hard to quantify. Low dimensional sufficient statistics would avoid this problem but generally do not exist, and there are many competing methods to choose summaries which make a good trade-off between low dimension and informativeness (Blum et al., 2013; Prangle, 2015). An alternative approach of Nott et al. (2014) is to improve ABC output by adjusting each parameter’s margin to agree with a separate marginal ABC analysis. These analyses can each use different low dimension summary statistics, so that the effect of the curse of dimensionality on the margins is reduced. However there are still issues in selecting these summaries,

and dealing with approximation error in the dependence structure.

Several other authors have recently investigated latent variable approaches to ABC. Neal (2012) introduced *coupled ABC* for household epidemics. This simulates latent variable vectors from their prior and, for each, finds one or many parameter vectors leading to closely matching simulated datasets. These parameters, weighted appropriately, form a sample from an approximate posterior. A similar strategy is employed for more general applications in the *reverse sampler* of Forneron and Ng (2015) and Meeds and Welling (2015)’s *optimisation Monte Carlo*. Alternatively, Moreno et al. (2016) perform variational inference, using latent variable vectors drawn from their prior in the estimation of loss function gradients.

A similar SMC approach to ours is outlined, but not implemented, by Andrieu et al. (2012). Analogous methods have been implemented for ABC inference of state space models, using ABC particle filtering to estimate likelihoods for a sequence of observations (Jasra, 2015). Another related method is Graham and Storkey (2016), who sample from the  $(\theta, x)$  space conditioned exactly on the observations using constrained Hamiltonian Monte Carlo (HMC). A limitation is that the  $y(\theta, x)$  mapping must be differentiable with respect to both arguments.

## 1.2 Contributions and overview

We provide an approximate inference method for the same class of intractable problems as ABC. Our algorithm samples from the same family of posterior approximations as ABC, but can reach more accurate approximations for the same computational cost. In particular, its cost rises more slowly with the data dimension. Therefore it is feasible to perform inference using a larger, and hence more informative, set of summary statistics. In some cases it is even feasible to use the full data.

Our method has various differences to competing methods using latent variables. Unlike the majority of these, it does not rely solely on randomly sampling latent variables, but instead searches their space more efficiently. Also unlike HMC approaches we do not require

differentiability assumptions for  $y(\theta, x)$ .

Typically SMC methods have many tuning choices. Another benefit of our approach is that these can all be automated. The tuning choices required are simply those for the ABC and PMMH algorithms.

Section 2 describes background information on the methods we use. Section 3 presents our algorithm to estimate the likelihood given a particular parameter vector, and how we use this within a MCMC inference algorithm. Asymptotic results on computational cost are also given here, quantifying the improvement over standard ABC. The method is evaluated on a simple Gaussian example in Section 4, and used in an infectious disease application in Section 5. Code for these examples is available at <https://github.com/dennisprangle/RareEventABC.jl>. Section 6 gives a concluding discussion, including when we expect our scheme to work well. Appendix A contains technical details of our asymptotics.

## 2 Background

### 2.1 Approximate Bayesian Computation

Suppose observations  $y_{\text{obs}}$  are available and we wish to learn the parameters  $\theta$  of a model  $\pi(y|\theta)$  (a density with respect to a probability measure  $dy$ ) given a prior  $\pi(\theta)$  (a density with respect to probability measure  $d\theta$ ). Algorithm 1 is an ABC rejection sampling algorithm which performs approximate Bayesian inference. It requires three tuning choices: the number of simulations  $N$ , a threshold  $\epsilon \geq 0$ , and a distance function  $d(\cdot, \cdot)$ . The latter is typically Euclidean distance or a variation. It is usually sensible to scale data  $y$  appropriately so that all components make contributions of similar size to the distance, and we will assume that this has already been done.

The output of Algorithm 1 is a sample from the following approximate posterior density

$$\pi_{\text{ABC}}(\theta|y_{\text{obs}}; \epsilon) \propto \pi(\theta)L_{\text{ABC}}(\theta; \epsilon), \quad (1)$$

---

**Algorithm 1** ABC rejection sampling

---

**Loop over**  $i = 1, 2, \dots, N$ .

1. Sample  $\theta_i$  from  $\pi(\theta)$ .
2. Sample  $y_i$  from  $\pi(y|\theta_i)$ .
3. Accept if  $d(y_i, y_{\text{obs}}) \leq \epsilon$ .

**End loop**

4. **Return:** accepted  $\theta_i$  values.
- 

where

$$L_{\text{ABC}}(\theta; \epsilon) = V(\epsilon)^{-1} \int \pi(y|\theta) \mathbb{1}[d(y, y_{\text{obs}}) \leq \epsilon] dy, \quad (2)$$

$$V(\epsilon) = \int \mathbb{1}[d(y, y_{\text{obs}}) \leq \epsilon] dy. \quad (3)$$

The *ABC likelihood*  $L_{\text{ABC}}$  is a convolution of the exact likelihood function and the *kernel*

$$k(y; \epsilon) = V(\epsilon)^{-1} \mathbb{1}[d(y, y_{\text{obs}}) \leq \epsilon], \quad (4)$$

a uniform density on  $y$  values close to  $y_{\text{obs}}$ . Under some weak conditions, as  $\epsilon \rightarrow 0$  the ABC likelihood converges to the exact likelihood and  $\pi_{\text{ABC}}$  to the exact posterior (This is shown by equation (6) in Appendix A, which describes some sufficient conditions.) However this causes acceptances to become rare. Thus there is a trade-off, controlled by  $\epsilon$ , between output sample size and the accuracy of  $\pi_{\text{ABC}}$ .

ABC rejection sampling is inefficient in the common situation where the prior is much more diffuse than the posterior, as a lot of time is spent on simulations that have very little chance of being accepted. Several more sophisticated ABC algorithms have been proposed which concentrate on performing simulations for  $\theta$  values believed to have high posterior density. These include versions of importance sampling, MCMC and SMC. These also output samples (sometimes weighted) from an approximation to the posterior, usually  $\pi_{\text{ABC}}$  as given

in (1). See Marin et al. (2012) for a review of ABC, including these algorithms and related theory.

As mentioned earlier, ABC suffers from a *curse of dimensionality* issue. Intuitively, the problem is that simulations producing good matches of all summaries simultaneously become increasingly unlikely as  $\dim(y)$  grows. For Algorithm 1, it has been proved (Blum, 2010; Barber et al., 2015; Biau et al., 2015) that for a fixed value of  $N$  the quality of the output sample as an approximation of the posterior deteriorates as  $d$  increases, even taking into account the possibility of adjusting  $\epsilon$ . See Fearnhead and Prangle (2012) for heuristic arguments that the problem also applies to other ABC algorithms.

## 2.2 Rare event sequential Monte Carlo

C  rou et al. (2012) propose two algorithms for estimating rare event probabilities using a SMC approach. We state these as Algorithms 2 (RE-SMC) and 3 (adaptive RE-SMC), using notation which will be helpful later for use in an ABC setting. The aim is to estimate a small probability  $L = \Pr(\Phi(x) \leq \epsilon|\theta)$ . Here  $x$  is a random variable,  $\theta$  is a vector of parameters,  $\Phi$  maps  $x$  values to  $\mathbb{R}$ , and  $\epsilon$  is a threshold. In the ABC setting of later sections  $L$  will be related to likelihood estimation.

C  rou et al. (2012) prove that RE-SMC produces an unbiased estimator of  $L$ , but adaptive RE-SMC gives an estimator with  $O(N^{-1})$  bias. They also analyse the asymptotic variance of the estimator’s relative error for large  $N$  under various assumptions. This variance is generally smaller for the adaptive estimator. Equality occurs only when RE-SMC uses an  $\epsilon$  sequence such that  $\Pr(\Phi(x) \leq \epsilon_{k+1}|\theta, \Phi(x) \leq \epsilon_k)$  is constant. An approximation to this sequence can be generated by running adaptive RE-SMC. We discuss which RE-SMC algorithm to use within our method later. Under optimal conditions the relative error variances decrease with  $T$ , the number of iterations, so that the estimates are more accurate than using plain Monte Carlo, which corresponds to  $T = 1$ . This result could be extended to take computational cost into account. However instead we will analyse the overall efficiency of

our proposed approach in Section 3.3.

---

**Algorithm 2** Rare event SMC algorithm (RE-SMC)

---

**Input:** Parameters  $\theta$ , number of particles  $N$ , thresholds  $\epsilon_1, \epsilon_2, \dots, \epsilon_T$ , Markov kernels for step 3.

1. For  $i = 1, \dots, N$  sample  $x_0^{(i)}$  from  $\pi(x|\theta)$ .

**Loop over**  $1 \leq t \leq T$ :

2. Calculate  $I_t = \{i | \Phi(x_{t-1}^{(i)}) \leq \epsilon_t\}$ . Let  $\hat{L}_t = |I_t|/N$ .  
(If  $\hat{L}_t = 0$  terminate algorithm returning  $\hat{L} = 0$ .)
3. For  $i = 1, \dots, N$  sample  $x_t^{(i)}$  by drawing  $j$  uniformly from  $I_t$  and applying a Markov kernel to  $x_{t-1}^{(j)}$  with invariant density  $\pi(x|\theta, \Phi(x) \leq \epsilon_{t-1})$  (taking  $\epsilon_0 = \infty$ ).

**End loop**

4. **Return:**  $\hat{L} = \prod_{t=1}^T \hat{L}_t$ .
- 

**Remarks**

1. In adaptive RE-SMC, typically  $N_{\text{acc}}$  particles are accepted so that  $|I_t| = N_{\text{acc}}$ . However there may be more acceptances in the final iteration or if ties in distance are possible.
2. For  $t \leq T$ ,  $\prod_{\tau=1}^t \hat{L}_\tau$  is an upper bound on  $\hat{L}$  in RE-SMC. This bound can be calculated during the  $t$ th iteration of the algorithms. This will be used later to terminate the algorithms early once the estimate is guaranteed to be below some prespecified bound.
3. The  $x_T^{(i)}$  values can be used for inference of  $x|\theta, \Phi(x) \leq \epsilon$ . When this is not of interest, as in this paper, then the computational cost can be reduced by omitting step 3 (resampling and Markov kernel propagation) in the final iteration of either algorithm.
4. It's possible for adaptive RE-SMC not to terminate. This could occur if the  $x_t^{(i)}$  particles become stuck near a mode where  $\Phi(x) > \epsilon$  and the Markov kernel is unable to move them to other modes. In Section 3.2 we will discuss how our proposed method can



---

**Algorithm 3** Adaptive rare event SMC algorithm (adaptive RE-SMC)

---

**Input:** Parameters  $\theta$ , number of particles  $N$ , target number to accept  $N_{\text{acc}}$ , acceptance thresholds  $\epsilon$ , rule to generate Markov kernels for step 3.

1. For  $i = 1, \dots, N$  sample  $x_0^{(i)}$  from  $\pi(x|\theta)$ .

**Loop over**  $t = 1, 2, \dots$ :

2. Let  $\epsilon_t$  be the maximum of (a) the  $N_{\text{acc}}$ th smallest  $\Phi(x_{t-1}^{(i)})$  value and (b)  $\epsilon$ . Calculate  $I_t = \{i | \Phi(x_{t-1}^{(i)}) \leq \epsilon_t\}$  and  $\hat{L}_t = |I_t|/N$ .
3. For  $i = 1, \dots, N$  sample  $x_t^{(i)}$  by drawing  $j$  uniformly from  $I_t$  and applying a Markov kernel to  $x_{t-1}^{(j)}$  with invariant density  $\pi(x|\theta, \Phi(x) \leq \epsilon_{t-1})$  (taking  $\epsilon_0 = \infty$ ).
4. If  $\epsilon_t = \epsilon$  break loop and go to step 5, setting  $T = t$ .

**End loop**

5. **Return:**  $\hat{L} = \prod_{t=1}^T \hat{L}_t$ .
- 

avoid this problem by terminating once it becomes clear the final likelihood estimate will be very low.

5. When ties in the distance are possible, adaptive RE-SMC iterations can fail to reduce the threshold. That is, sometimes step 2 can give  $\epsilon_{t+1} = \epsilon_t$ . This can produce very long run times. Possible improvements to deal with this are discussed in Section 6. (Note that when adaptive RE-SMC is being used to select a sequence of thresholds then repeated values should be removed.)
6. These algorithms use multinomial resampling. More efficient schemes exist, but are not investigated by the theoretical results of Cérou et al. (2012).

## 2.3 Slice sampling

We require a suitable Markov kernel to use within the RE-SMC algorithms. This must have invariant density  $\pi(x|\theta, \Phi(x) \leq \epsilon_{t-1})$ . As discussed below in Section 3.2, our ABC setting will assume  $\pi(x|\theta)$  is uniform on  $[0, 1]^m$ . Hence the required invariant distribution is uniform on

the subset of  $[0, 1]^m$  such that  $\Phi(x) \leq \epsilon_{t-1}$ . We will use *slice sampling* as the Markov kernel. This section outlines the general idea of slice sampling, then how it can be implemented in this setting and its advantages over alternative choices.

Slice sampling is a family of MCMC methods to sample from an unnormalised target density  $\gamma(x)$ . The general idea is to sample uniformly from the set  $\{(x, h) \mid h \leq \gamma(x)\}$  and marginalise. We will concentrate on an algorithm of Murray and Graham (2016) for the case where the support of  $\gamma(x)$  is  $[0, 1]^m$ , or a subset of this. Their algorithm updates the current state  $x$  by first drawing  $h$  from  $\text{Uniform}(0, \gamma(x))$ , then proposing  $x'$  values, accepting the first one for which  $\gamma(x') \geq h$ . The proposal scheme initially considers large changes from  $x$  in a randomly chosen direction, and then, if these are rejected, progressively smaller changes.

For use within RE-SMC,  $\gamma(x)$  can be taken to be the indicator function  $\mathbb{1}(\Phi(x) \leq \epsilon_{t-1})$ . This means the condition  $\gamma(x') \geq h$  simplifies to  $\gamma(x') > 0$ , so sampling  $h$  can be omitted. The resulting slice sampling update is given by Algorithm 4, which is a special case of the Murray and Graham (2016) algorithm mentioned above (and similar to the *hit-and-run sampler*; see Smith, 1996). See their paper for details of the proof that  $\gamma(x)$  is the invariant density of this Markov kernel.

Next we describe two advantages of using slice sampling within RE-SMC, particularly in relation to the alternative of using a Metropolis-Hastings kernel. Firstly, slice sampling requires little tuning. If tuning choices were required, for example a proposal distribution for Metropolis-Hastings, then RE-SMC would need to include rules to make a good choice automatically, which may be difficult. Another advantage of slice sampling is that each iteration outputs a unique  $x$  value. On the other hand, Metropolis-Hastings can produce duplicate values, which is problematic within SMC because it leads to increased variance of probability estimates.

A tuning choice which is required by slice sampling is the initial search width  $w$ . A default choice is  $w = 1$ , but this means that the number of loops required will increase for small  $\epsilon$ . To deal with this we choose  $w = 1$  in the first SMC iteration and then select  $w$

adaptively, as  $\min(1, 2\bar{z})$  where  $\bar{z}$  is the maximum final value of  $|z|$  from all slice sampling calls in the previous SMC iteration. This choice generally shrinks  $w$  based on the most recent value of  $\bar{z}$ , while avoiding some unwanted behaviours. Firstly it avoids forcing  $w$  to decrease at a fixed rate, so that eventually only very small steps would be attempted. Secondly it avoids  $w$  growing above 1, which would make slice sampling expensive when local moves are required. The effect of our choice is investigated empirically later (see Figure 3).

---

**Algorithm 4** Slice sampling update for RE-SMC

---

**Input:** current state  $x$  of dimension  $p$ , map  $\Phi(x)$ , threshold  $\epsilon$ , initial search width  $w$ . It's assumed that  $\Phi(x) \leq \epsilon$ .

1. Sample  $v \sim N(0, I_p)$
2. Sample  $u \sim \text{Uniform}(0, w)$ . Let  $a = -u, b = w - u$ .

**Loop:**

3. Sample  $z \sim \text{Uniform}(a, b)$ .
4. Define a vector  $x'$  by  $x'_i = r(x_i + zv_i)$  using the *reflection function*:

$$r(y) = \begin{cases} m & m < 1 \\ 2 - m & m \geq 1 \end{cases}$$

where  $m$  is the remainder of  $y$  modulo 2.

5. If  $\Phi(x') \leq \epsilon$  then **return**  $x'$ .
6. If  $z < 0$  let  $a = z$ , otherwise let  $b = z$ .

**End loop**

---

## 2.4 Pseudo-marginal Metropolis-Hastings

The Metropolis-Hastings (MH) algorithm samples from a Markov chain with stationary distribution proportional to an unnormalised density  $\psi(\theta)$ . It is often used in Bayesian inference to produce samples from a close approximation to the posterior distribution. Despite the non-independence of these samples, they can still be used to produce highly accurate Monte

Carlo estimates of functions of the posterior. Simulating  $\theta_t$ , the  $t$ th state of the Markov chain is based on sampling a state  $\theta'$  from a proposal density  $q(\theta'|\theta_{t-1})$ , typically centred on the preceding state  $\theta_{t-1}$ . This proposal is accepted as  $\theta_t$  with probability  $\min\left(1, \frac{\psi(\theta')q(\theta_{t-1}|\theta')}{\psi(\theta_t)q(\theta'|\theta_{t-1})}\right)$ . Otherwise  $\theta_t = \theta_{t-1}$ .

This algorithm remains valid if likelihood evaluations are replaced with unbiased non-negative estimates as follows (Andrieu and Roberts, 2009). The state of the Markov chain is now  $(\theta_t, \hat{\psi}_t)$ , where  $\hat{\psi}_t$  is an estimate of  $\psi(\theta_t)$ , and the acceptance probability must be  $\min\left(1, \frac{\hat{\psi}'q(\theta_{t-1}|\theta')}{\hat{\psi}_{t-1}q(\theta'|\theta_{t-1})}\right)$ . Crucially, upon acceptance  $\hat{\psi}_t$  is set to the estimate  $\hat{\psi}'$  for the proposal  $\theta'$ . So, rather than being recalculated in every iteration, this estimate is used in all future iterations until another proposal is accepted. A version of the resulting *pseudo-marginal Metropolis-Hastings* (PMMH) algorithm, specialised to this paper’s setting, is presented below as Algorithm 5.

Optimal tuning of PMMH has been examined theoretically by Pitt et al. (2012), Doucet et al. (2015) and Sherlock et al. (2015), covering the case where each  $\hat{\psi}'$  estimate is generated by an SMC algorithm. A central issue is how many SMC particles should be used to optimise the computational efficiency of PMMH. All the authors conclude that this number should be tuned to achieve a particular variance of  $\log \hat{\psi}$ . (It’s assumed, unrealistically, that this variance does not depend on  $\theta$ . In practice it’s typical to investigate the variance at a fixed value of  $\theta$  believed to have high posterior density.) The value derived for this optimal variance differs between the authors due to their different assumptions, but all values lie in the range 0.8–3.3. Sherlock et al. (2015) also investigate tuning the proposal distribution  $q$ , and suggest using proposal variance  $\frac{2.562^2}{\dim(\theta)}\Sigma$  where  $\Sigma$  is the posterior variance. They perform simulation studies generally supporting both these results. One key assumption made by all the authors is that  $\log \hat{\psi}$  follows a normal distribution. The validity of this assumption in our setting will be investigated later. It’s also assumed that the computational cost of SMC is proportional to the number of particles used and does not depend on  $\theta$ , which is generally true for SMC algorithms.

### 3 High dimensional ABC

This section presents our approach to inference in the ABC setting, using the algorithms reviewed in Section 2. Section 3.1 describes how the RE-SMC algorithms can estimate the ABC likelihood given values of  $\theta$  and  $\epsilon$ . Such likelihood estimators can be used within several inference algorithms to produce approximate Bayesian inference. In this paper we concentrate on PMMH. Section 3.2 presents the resulting method. Section 3.3 discusses the computational cost of the resulting *RE-ABC algorithm* in comparison to standard ABC, with particular note of the high dimensional case.

#### 3.1 Likelihood estimation

For now, suppose  $\theta$  and  $\epsilon > 0$  are fixed. We aim to produce an unbiased estimate of  $L_{\text{ABC}}(\theta; \epsilon)$ , as defined in (2).

Suppose there exist latent variables  $x$  such that the observations can be written as a deterministic function  $y = y(\theta, x)$ . The idea is that  $x$  and  $\theta$  suffice to specify a complete realisation of the simulation process, even including details such as observation error, and  $y(\theta, x)$  is a vector of partial observations. Neglecting  $\theta$ , which is fixed for now,  $y(\theta, x)$  will be written below as simply  $y = y(x)$ . See Section 6 for a discussion of properties of  $y(x)$  which help our approach work well.

We specify a density  $\pi(x|\theta)$  (with respect to Lebesgue measure) for the latent variables. This is part of the specification of the model, but it can also be viewed as representing prior beliefs about the latent variables. Throughout the paper we take  $\pi(x|\theta)$  to be uniform on  $[0, 1]^m$  regardless of  $\theta$ . Under this interpretation,  $x$  is a vector of  $m$  independent standard uniform random variables which suffice to carry out the simulation process.

Now we simply apply RE-SMC or adaptive RE-SMC using  $\Phi(x) = d(y(x), y_{\text{obs}})$ . The

small probability estimated by these algorithms is

$$\begin{aligned}\Pr(\Phi(x) \leq \epsilon|\theta) &= \int \pi(x|\theta) \mathbb{1}[d(y(x), y_{\text{obs}}) \leq \epsilon_t] dx \\ &= \int \pi(y|\theta) \mathbb{1}[d(y, y_{\text{obs}}) \leq \epsilon_t] dy,\end{aligned}$$

which equals  $L_{\text{ABC}}(\theta; \epsilon)$  multiplied by the constant  $V(\epsilon)$ . Hence using RE-SMC we can obtain an estimate of  $L_{\text{ABC}}(\theta; \epsilon)$  which is unbiased, as required by PMMH. Using adaptive RE-SMC produces a slightly biased estimate, and we comment on the effect of using this within PMMH in the next section.

Note that for  $\pi(x|\theta)$  uniform, the implementation of slice sampling in Algorithm 4 is well-suited to be the Markov kernel within RE-SMC. For other  $\pi(x|\theta)$  distributions, alternative Markov kernels would usually be necessary, for example elliptical slice sampling (see Murray and Graham, 2016) for the Gaussian case, or Gibbs updates for the discrete case.

### 3.2 Inference

Algorithm 5 shows the PMMH algorithm for our setting, which we refer to as RE-ABC. The probability of acceptance in step 4 corresponds to a target density proportional to  $\pi(\theta)L_{\text{ABC}}(\theta; \epsilon)$  i.e. the ABC posterior (1).

The user must decide whether to use RE-SMC or adaptive RE-SMC within RE-ABC to estimate the ABC likelihood (Later we sometimes refer to “adaptive” and “non-adaptive” RE-ABC depending on which is used.) RE-SMC is unbiased so will more faithfully reproduce the results of ABC. However the bias introduced is small, and may have little effect compared to the efficiency benefits of the variance reduction which adaptive RE-SMC provides (theoretical and practical aspects of MCMC algorithms that have this character are discussed in Alquier et al. 2016). We investigate this empirically in Sections 4 and 5 and find no noticeable effect of bias. However we find adaptive RE-SMC to sometimes be less computationally efficient in practice, and so we recommend using the non-adaptive algorithm.

---

**Algorithm 5** Pseudo-marginal Metropolis-Hastings using RE-SMC (RE-ABC)

---

**Input:** initial state  $\theta_1$ , number of iterations  $M$ , number of SMC particles  $N$  and tuning choices for PMMH and ABC.

1. Let  $t = 1$  and calculate  $\hat{L}_1$  using RE-SMC (or adaptive RE-SMC) with slice sampling as the Markov kernel.

**Loop over**  $2 \leq t \leq M$ :

2. Propose new state  $\theta'$  from  $q(\cdot|\theta)$  and sample  $u$  from Uniform(0, 1).
3. Calculate  $\hat{L}'$  using RE-SMC (or adaptive RE-SMC) with slice sampling as the Markov kernel. This calculation can be stopped early once rejection in the next step is guaranteed.
4. If  $u > \frac{\pi(\theta')\hat{L}'q(\theta_{t-1}|\theta')}{\pi(\theta_{t-1})\hat{L}_{t-1}q(\theta'|\theta_{t-1})}$ :

*Reject* Let  $\theta_t = \theta_{t-1}$  and  $\hat{L}_t = \hat{L}_{t-1}$ .

Else:

*Accept* Let  $\theta_t = \theta'$  and  $\hat{L}_t = \hat{L}'$ .

**End loop**

**Return:**  $\theta_1, \theta_2, \dots, \theta_M$ .

---

Reasons for this are described shortly, and discussed in more detail, along with possibilities for improvement, in Section 6.

To reduce computational costs RE-SMC (adaptive or non-adaptive) can be terminated as soon as rejection is guaranteed. To implement this, after step 2 check whether

$$\prod_{\tau=1}^{t_{\text{SMC}}} \hat{L}_{\tau} < \frac{u\pi(\theta_{t-1})\hat{L}_{t-1}q(\theta'|\theta_{t-1})}{\pi(\theta')q(\theta_{t-1}|\theta')},$$

where  $t_{\text{SMC}}$  is the value of the  $t$  variable within the RE-SMC algorithm. If this is true, terminate the RE-SMC algorithm and reject the current proposal in the PMMH algorithm. The MCMC algorithm remains valid since the final RE-SMC likelihood estimate is guaranteed to be smaller than  $\prod_{\tau=1}^{t_{\text{SMC}}} \hat{L}_{\tau}$  and therefore lead to rejection in PMMH. Early termination prevents extremely long runs of RE-SMC for  $\theta$  values with low posterior densities. It is most efficient for non-adaptive RE-SMC, where it is always possible to terminate in any iteration if the  $\hat{L}_{\tau}$  values are small enough. For adaptive RE-SMC,  $\hat{L}_{\tau} \geq \frac{N_{\text{acc}}}{N}$  so there is a lower bound of how many iterations are required before termination. This argument suggest adaptive RE-SMC is less computationally efficient, and agrees with later empirical findings (see Figure 5).

Earlier we commented that adaptive RE-SMC can fail to terminate in some situations. When ties in the distance are not possible, then this is usually not a problem within RE-ABC due to the early termination rule just outlined. However care is still required the first time adaptive RE-SMC is run, and when it is used in pilot runs. Ties in the distance can be more problematic and are discussed further in Section 6.

There are numerous tuning choices required in this PMMH algorithm. Most of these can be based on the output of a pilot analysis, for example an ABC analysis or a short initial run of PMMH. The estimated posterior mean  $\hat{\mu}$  can be used as an initial PMMH state. The estimated posterior variance  $\hat{\Sigma}$  can be used to tune the PMMH proposal density. Following the PMMH theory discussed in Section 2.4 we sample proposal increments from



$N\left(0, \frac{2.562^2}{\dim(\theta)} \hat{\Sigma}\right)$ . (Note that the early termination rule avoids SMC calls having very long run times for some  $\theta$  values, approximately meeting the assumptions of the PMMH tuning literature.) The threshold sequence for RE-SMC can be selected by running adaptive RE-SMC with  $\theta = \hat{\mu}$ . To select the number of particles, a few preliminary runs of RE-SMC (or adaptive RE-SMC) can be performed with  $\theta = \hat{\mu}$ , aiming to produce a log likelihood variance of roughly 1. This is at the more conservative end of the range suggested by the theory reviewed earlier.

A crucial tuning choice which remains is  $\epsilon$ . As in other ABC methods, we suggest tuning this pragmatically based on the computational resources available. This can be done by running adaptive RE-SMC with  $\theta = \hat{\mu}$  and  $\epsilon = 0$  and stopping after a prespecified time, corresponding to how long is available for an iteration of PMMH. The value of  $\epsilon_t$  when the algorithm is stopped can be used as  $\epsilon$ . It is still possible for the SMC algorithms to take much longer to run for other  $\theta$  values. However the early termination rule will usually mitigate this. Diagnostic plots can be used to investigate whether the  $\epsilon$  value selected produces simulations judged to be sufficiently similar to the observations. For example, see Figure 1 of the supplementary material.

### 3.3 Cost

Here we summarise results on the cost of ABC and RE-ABC in terms of time per samples produced (or effective sample size for PMMH algorithms), in the asymptotic case of small  $\epsilon$ . Arguments supporting these results are given in Appendix A. Several assumptions are required, principally that the density  $\pi(y|\theta)$  is with respect to Lebesgue measure – informally, the observations must be continuous. Weakening these assumptions is discussed in supplementary material. Note that the results for RE-ABC are the same whether RE-SMC or adaptive RE-SMC is used.

The time per sample is asymptotic to  $1/V(\epsilon)$  for ABC and  $[\log V(\epsilon)]^2$  for RE-ABC (see (3) for definition of  $V(\epsilon)$ .) So, asymptotically, RE-ABC has a significantly lower cost to reach the

same target density. To illustrate the effect of  $D = \dim(y)$  we can consider the asymptotic case of large  $D$  (n.b. as shown in the supplementary material, when some observations are non-continuous then  $D$  can be replaced with the dimension of  $\{y | d(y, y_{\text{obs}}) < \epsilon\}$  for small  $\epsilon$ .) Under the Lebesgue assumption, (3) gives that  $V(\epsilon) \propto \epsilon^D$ . Hence the time per sample is asymptotic to the following expressions, written in terms of  $\tau = 1/\epsilon$  for interpretability:  $C_1 = \tau^D$  for ABC and  $C_2 = D^2[\log \tau]^2 = [\log C_1]^2$  for RE-ABC. Hence ABC has an exponential cost in  $D$ , while RE-ABC has only a quadratic cost. This makes high-dimensional inference more tractable for RE-ABC but dimension reduction via summary statistics will remain useful in controlling the cost when  $D$  is large.

These results assume the algorithms are run sequentially. The PMMH stage of RE-ABC is innately sequential, but particle updates can be run in parallel, providing a benefit from parallelisation. Compared to the most efficient ABC algorithms, this is an advantage over ABC-MCMC and seems roughly comparable to that of ABC-SMC algorithms.

## 4 Gaussian example

In this section we compare ABC (Algorithm 1) and RE-ABC (Algorithm 5) on a simple Gaussian model. The model is  $Y_i \sim N(0, \sigma^2)$  independently for  $1 \leq i \leq 25$ . We use the prior  $\sigma \sim \text{Uniform}(0, 10)$ . This is an interesting test case because  $\dim(y)$  is large enough to cause difficulties for ABC methods but calculations are quick, and the results can be compared to those of likelihood-based methods.

### 4.1 Comparison of ABC and RE-ABC

We compared ABC and RE-ABC for observations drawn from the model using  $\sigma = 3$ . For each of  $\epsilon = 8, 9, \dots, 30$ , we ran ABC until  $N = 500$  simulations were accepted and calculated the root mean squared error and time per acceptance. Both adaptive and non-adaptive versions of RE-ABC were run for 2000 iterations with  $\epsilon = 5, 10, 15, 20, 25$ . As described in

Section 3.2, pilot runs were used to tune the number of particles, the Metropolis-Hastings proposal standard deviation and, where necessary, the threshold sequence. We chose the number of acceptances in all adaptive RE-ABC analyses to be half the number of particles. To avoid dealing with burn-in, we started the PMMH chains at  $\sigma = 3$ . For comparison we also ran ABC-MCMC (Marjoram et al., 2003) and MCMC using the exact likelihood.

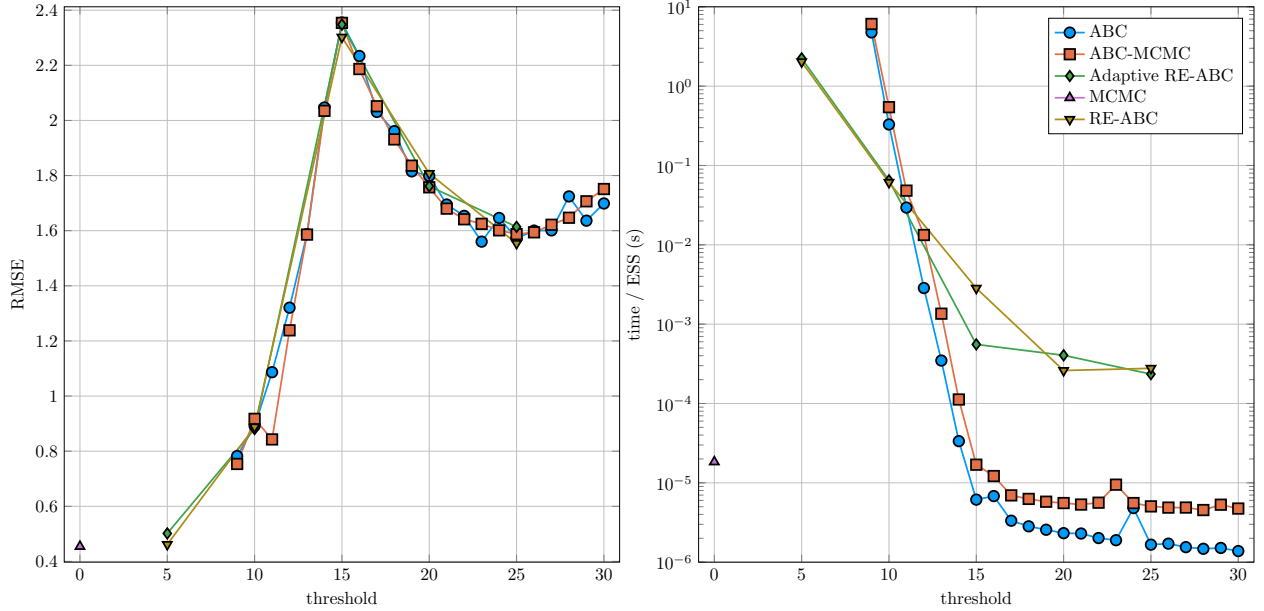


Figure 1: Simulation study comparing ABC, ABC-MCMC, adaptive and non-adaptive RE-ABC, and exact likelihood MCMC on IID Gaussian data.

Figure 1 shows the results. The left panel illustrates that accuracy improves as the acceptance threshold  $\epsilon$  is reduced below roughly 15, and, as expected, all methods produce very similar results. The right panel investigates the time taken per sample by ABC. For MCMC output, this is time divided by the effective sample size (the IMSE estimate of Geyer, 1992.) Under ABC and ABC-MCMC, time per sample increases rapidly as  $\epsilon$  is reduced. For both RE-ABC algorithms the increase is slower, allowing smaller values of  $\epsilon$  to be investigated. Neither RE-ABC algorithm is obviously more efficient than the other. This difference between ABC and RE-ABC is consistent with the asymptotics on computational cost described in Section 3.3. However for large  $\epsilon$  values ABC and ABC-MCMC are cheaper.

Overall RE-ABC permits smaller  $\epsilon$  values to be investigated at a reasonable computational cost, producing more accurate approximations.

Figure 2 provides some further insight into the efficiency of the RE-ABC algorithms, by looking at the times taken for calls to the RE-SMC algorithm. These have similar distributions for the adaptive and non-adaptive algorithms, indicating that there is little difference in their efficiency. One point of interest is that the adaptive algorithm takes a minimum time of 0.095 seconds even when it stops early, while the non-adaptive algorithm sometimes stops early in a much shorter time. However this happens too rarely to have much effect on overall efficiency.

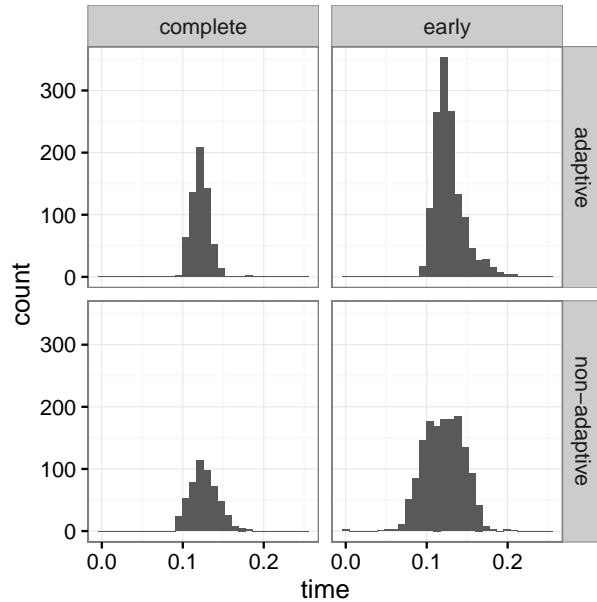


Figure 2: Histograms of times (in seconds) taken by calls to RE-SMC within an RE-ABC analysis of IID Gaussian data. The top row is for an analysis using the adaptive algorithm, and the bottom for the non-adaptive algorithm. Both analyses used  $\epsilon = 5$  and the same tuning details, chosen using a pilot run. The left column is for those calls in which RE-SMC was completed, while the right shows those where it was terminated early.

## 4.2 Validity of assumptions

We also used the Gaussian example to investigate the validity of various assumptions about RE-ABC used in this paper. First we considered the cost of slice sampling calls in RE-ABC. Figure 3 shows the mean number of iterations that slice sampling requires during an illustrative RE-ABC run, with a fixed  $\epsilon$  sequence. Two cases are shown: non-adaptive slice sampling tuning ( $w = 1$  in Algorithm 4) or adaptive tuning ( $w$  updated as described in Section 2.3). This gives empirical evidence that adaptive tuning prevents the slice sampling cost from increasing during the algorithm, as desired. Repeated trials show that both methods produce very similar mean likelihoods. However adaptive tuning did increase the log-likelihood variance slightly so there is a small trade-off in its use.

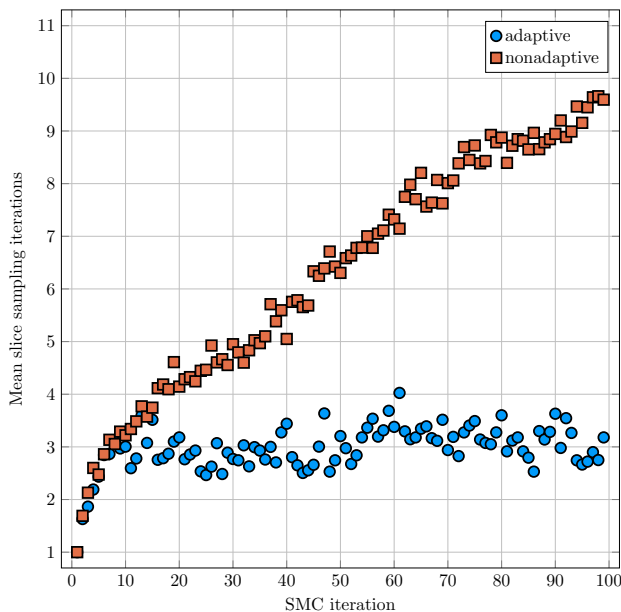


Figure 3: Number of slice sampling iterations required under adaptive and non-adaptive rules for selecting the tuning parameter  $w$  within an RE-SMC run on IID Gaussian data..

Secondly we investigated the distribution of likelihood estimates produced by RE-ABC. Recall that the theoretical literature on PMMH assumes that these follow a log-normal distribution. Figure 4 shows quantile-quantile plots comparing log likelihood estimates to normal quantiles. The estimates are approximately normal when a sufficient number of

particles are used, but become increasingly skewed as this shrinks. A major departure from normality is that for a small number of particles many likelihood estimates are zero. The corresponding points are omitted from the plot. In conclusion, the normality assumption seems reasonable if a sufficient number of particles are used.

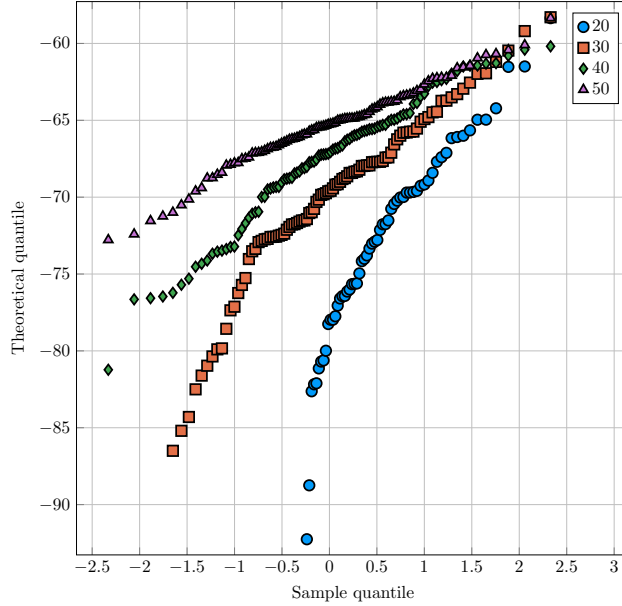


Figure 4: Normal quantile-quantile plots of log likelihood estimates for an RE-ABC example under various numbers of particles. Omitted points correspond to likelihood estimates of zero.

## 5 Epidemic application

Infectious disease data is often modelled using *compartment models* where members of a population passes through several stages. We will consider a model with susceptible, infectious and removed stages – the so-called *SIR model* (Andersson and Britton, 2000). A susceptible individual has not yet been infected with the disease but is vulnerable. An infectious individual has been infected and may spread the disease to others. A removed individual can no longer spread the disease. Depending on the disease this may be due to immunity following recovery, or death.

We will use a stochastic version of this model based on a continuous time stochastic process  $\{S(t), I(t) : t \geq 0\}$  for numbers susceptible and infectious at time  $t$ . The total population size is fixed at  $n$  so the number removed at time  $t$  can be derived as  $R(t) = n - S(t) - I(t)$ . The initial conditions are  $(S(0), I(0)) = (n - 1, 1)$ . Two jump transitions are possible: infection  $(i, j) \mapsto (i - 1, j + 1)$  and removal  $(i, j) \mapsto (i, j - 1)$ . The simplest version of the model is Markovian and is defined by the instantaneous hazard functions of the two transitions, which are  $\frac{\lambda}{n}S(t)I(t)$  for infection and  $\gamma I(t)$  for removal. The unknown parameters are  $\lambda$ , controlling infection rates and  $\gamma$ , the removal rate. A goal of inference is often to learn about the basic reproduction number  $R_0 = \lambda/\gamma$ . This is the expected number of further infections caused by an initial infected individual in a large susceptible population. When  $R_0 < 1$ , most epidemics will infect an insignificant proportion of a large population. Many variations on the Markovian SIR model are possible, some of which are outlined below.

Likelihood-based inference is straightforward for fully observed data from an SIR model. However in practice only partial and possibly noisy observations of removal times are available, producing an intractable likelihood. For many models near-exact inference is possible by MCMC methods (summarised by McKinley et al., 2014), but small changes to the details require new and model-specific algorithms. Approximate inference can be performed by ABC (summarised by Kypraios et al., 2016), which is more adaptable but does not scale well to high-dimensional data. Here we illustrate how the RE-ABC algorithm can, without modification, perform inference for several variations on the SIR model, and do so more efficiently than standard ABC methods. As we concentrate on a classic and well-studied dataset, our analysis does not provide any novel subject-area insights.

Section 5.1 describes a method of simulating from SIR models. Section 5.2 discusses the distance function we use to implement RE-ABC. Data analysis is performed in Section 5.3.

## 5.1 Sellke construction

The Sellke construction (Sellke, 1983) for an SIR model provides an appealing way to simulate epidemic models. It introduces latent *infectious periods*  $g_i \sim F_{\text{inf}}$  and *pressure thresholds*  $p_i \sim F_{\text{press}}$  for  $1 \leq i \leq n$ , all independent. For the Markovian SIR model,  $F_{\text{inf}}$  is  $\text{Exp}(\gamma)$  and  $F_{\text{press}}$  is  $\text{Exp}(1)$ , but other choices are possible and may be more biologically plausible. We condition on  $g_1 = 0$  so that the first infection occurs at time 0. Algorithm 6 shows how these variables and the parameter  $\lambda$  are converted to simulated removal times. To use slice sampling we require the latent variables to be uniformly distributed a priori. Therefore we use quantiles of the  $g_i$ s and  $p_i$ s as the latent variables.

The cost of Algorithm 6 is  $O(n \log n)$ , where  $n$  is the population size. This is because the main loop runs at most  $2n - 1$  times, and involves finding the minimum of a set of up to  $n - 1$  removal times, which requires  $O(\log n)$  steps. (This is the case if the set is stored as an ordered vector. The cost of adding a new item is  $O(\log n)$ .)

Alternative simulation methods exist, principally the Gillespie algorithm (described in Kypraios et al., 2016, for example). Here the latent variables form a sequence controlling the behaviour of each successive jump event. The Gillespie algorithm has the advantage of  $O(n)$  cost. However it seems hard for slice sampling to explore the space of latent variables due to the behaviour of the mapping  $y(\theta, x)$ . In particular a small change in latent variables which alters the type of one jump will typically have a large and unpredictable effect on all the subsequent jumps. For more discussion on desirable properties of  $y(\theta, x)$ , see Section 6.

Note that when  $F_{\text{press}}$  is  $\text{Exp}(1)$  then  $R_0 = \lambda E(F_{\text{inf}})$  (Andersson and Britton, 2000). However to our knowledge the definition of  $R_0$  has not been extended to cover general  $F_{\text{press}}$ .

## 5.2 Distance function

Recall that the observations are the inter-removal times, or equivalently the times since the first removal. Let  $r_{(1)} \leq r_{(2)} \leq \dots \leq r_{(\nu)}$  denote the ordered removal times of a dataset with  $\nu$  removals. The times since first removal are then  $s_{(i)} = r_{(i)} - r_{(1)}$  for  $1 \leq i \leq \nu$ . We define



---

**Algorithm 6** Sellke construction epidemic simulator

---

Input: population size  $n$ , scaled infection rate parameter  $\beta = \lambda/n$ , infectious periods  $g_1, \dots, g_n$  and pressure thresholds  $p_2, \dots, p_n$ .

1. Set  $r_1 \mapsto g_1$  (assumes individual 1 has infection time 0).
2. Set  $r_i \mapsto \infty$  for  $i > 2$ .
3. Set  $I \mapsto 1$  (current number infected),  $t \mapsto 0$  (current time),  $p \mapsto 0$  (current pressure).
4. **While**  $I > 0$ :
  5. Find  $p_a = \min\{p_i | p_i > p\}$ . If this set is empty use  $p_{n+1} = \infty$ .
  6. Find  $r_b = \min\{r_i | r_i > t\}$ .
  7. Set  $p' \mapsto p + \beta I(r_b - t)$  (pressure at time  $r_b$  if  $I$  does not change)
  8. If  $p_a < p'$ :
    - (a) Set  $I \mapsto I + 1$ ,  $t \mapsto t + \frac{p_a - p}{\beta I}$ ,  $r_a \mapsto t + g_a$ ,  $p \mapsto p_a$ .
  9. Else:
    - (a) Set  $I \mapsto I - 1$ ,  $t \mapsto r_b$ ,  $p \mapsto p'$ .
10. **End while**

Output: Removal times  $r_1, r_2, \dots, r_n$ . Infinite removal time represents an individual who is never infected.

---

the distance between a simulated and observed dataset as:

$$\left[ \sum_{i \leq \min(\nu, \nu')} (s_{(i)} - s'_{(i)})^2 \right]^{1/2} + \sum_{\nu' < i \leq \nu'} [k + \bar{\rho}' - \rho'_{(i)}] + \sum_{\nu' < i \leq \nu} [k + \rho'_{(i)}] \quad (5)$$

where a dash denotes the simulated dataset. Here  $k$  is a tuning parameter penalising mismatches between  $\nu$  and  $\nu'$ . We take  $k = 1000$ . The  $\rho'_{(i)}$  terms are the sorted simulated pressure thresholds and  $\bar{\rho}'$  is the total simulated pressure (which equals  $\beta$  times the sum of the infectious periods for removed individuals). They are included to encourage these pressures to increase or decreasing appropriately to match  $\nu$  and  $\nu'$ . Without the pressure terms the adaptive version of RE-SMC sometimes failed to terminate. See Section 6 for further discussion.

### 5.3 Analysis of Abakaliki data

The Abakaliki dataset contains times between removals from a smallpox epidemic in which 30 individuals were infected from a closed population of 120. It has been studied by many authors under many variations to the basic SIR model. We study three models. The first model uses a  $\text{Gamma}(k, \gamma)$  infectious period (similar to Neal and Roberts, 2005). The second assumes pressure thresholds are distributed by a  $\text{Weibull}(k, 1)$  distribution (as in Streftaris and Gibson, 2012.) The third is the Markovian SIR model, but with removal times only recorded within 5 day bins. This is realised by altering the  $s_{(i)} - s'_{(i)}$  term (difference between simulated and observed day of removal) in (5) to  $f(s_{(i)}) - f(s'_{(i)})$  where  $f(s) = 5 \lfloor s/5 \rfloor$ , the greatest multiple of 5 less than or equal to  $s$ . In each model there are two or three unknown parameters:  $\lambda$ , controlling infection rates;  $\gamma$ , infectious period scale;  $k$ , a shape parameter. These are all assigned independent exponential prior distributions with rate 0.1, representing

weakly informative prior beliefs that these parameters are less likely to be large.

We chose the acceptance threshold to be  $\epsilon = 15$  on the pragmatic grounds that this produced run-times of no more than 6 hours on a desktop PC. Tuning was performed using pilot runs as described in Section 3.2. Of particular note is the number of particles required: 300 (Gamma infectious period), 200 (Weibull pressure thresholds) and 400 (binned removal times). Table 1 summarises the approximate posterior results. As the parameters differ between models, we don’t present parameter estimates. Instead we give several quantities of interest for each: the  $R_0$  estimate (where defined) and the means and standard deviations of (a) the pressure thresholds and (b) the infectious period. Most quantities are similar to each other and previous analyses (see McKinley et al., 2014 for a summary of many of these) despite the different modelling assumptions. A noticeable difference is that the infectious period is less variable in the model where it follows a Gamma distribution.

Figure 1 of the supplementary material shows simulated epidemics from each model. This shows that our choice of  $\epsilon$  produces epidemics reasonably close to the observed data for every model. Formal model choice is not straightforward in our framework (see discussion in Section 6), but it is easy to explore whether the models produced large differences in log-likelihood. In this case differences were modest, as shown by Figure 2 in the supplementary material, and within what would be explained, using BIC type arguments, by the differing number of parameters in the models. So we conclude qualitatively that there are no clear differences in fit between the models.

Model	$R_0$	Pressure thresholds		Infectious period	
		Mean	Standard deviation	Mean	Standard deviation
5 day bins	1.16 (0.30)	0.11 (0.03)	0.11 (0.03)	11.1 (3.0)	11.1 (3.0)
Gamma infectious period	1.18 (0.24)	0.09 (0.03)	0.09 (0.03)	13.6 (3.8)	6.8 (2.2)
Weibull pressure thresholds	–	0.10 (0.04)	0.11 (0.03)	12.4 (3.3)	12.4 (3.3)

Table 1: Approximate posterior estimates of basic reproduction number  $R_0$  and the means and standard deviations of pressure thresholds and infectious periods for the Abakaliki data under three models computing using RE-ABC. The table contains Monte Carlo estimates along with standard deviations in brackets. The  $R_0$  value is not given for the Weibull pressure threshold model as no definition is available for this model.

Adaptive RE-ABC was also tried and returned parameter inference results extremely similar to those for the non-adaptive algorithm – see Table 1 in the supplementary material. However, for some analyses the run times were longer. For example, the Gamma infectious period model took 263 minutes for the non-adaptive algorithm and 323 minutes for the adaptive algorithm. Figure 5 investigates this in more detail. It shows that run time difference is because most calls to RE-SMC terminate early, and these are generally quicker under the non-adaptive algorithm. It is also interesting that the adaptive algorithm is typically faster for completed RE-SMC calls. These findings are discussed in the next section.

We also ran ABC-MCMC for comparison, using the same MCMC and  $\epsilon$  tuning choices as for RE-ABC. The time per number of acceptance was at least 7 minutes in all cases. For RE-ABC this value was always less than 2 minutes. (Effective sample size was not used as there were sometimes too few ABC-MCMC acceptances to calculate it accurately.)

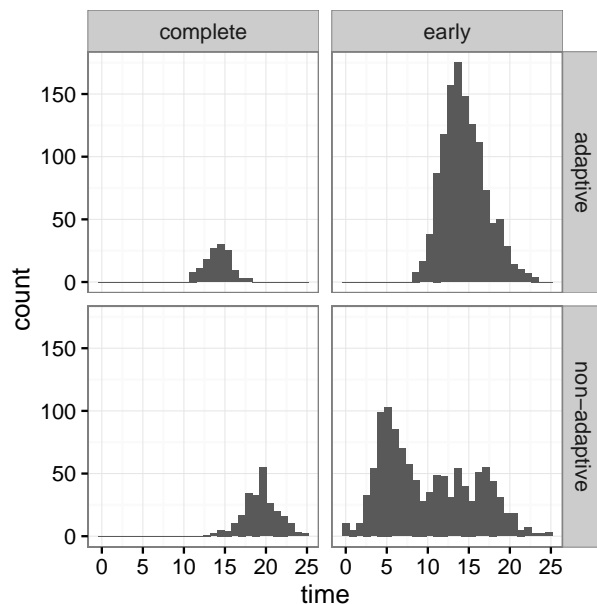


Figure 5: Histograms of times (in seconds) taken by calls to RE-SMC within an RE-ABC analysis of Abakaliki data. The top row is for an analysis using the adaptive algorithm, and the bottom for the non-adaptive algorithm. Both analyses used  $\epsilon = 15$  and the same tuning details, chosen using a pilot run. The left column is for those calls in which RE-SMC was completed, while the right shows those where it was terminated early.

## 6 Discussion

We have presented a method for approximate inference under an intractable likelihood when simulation of data is possible. It uses the same posterior approximation as ABC, (1), which is controlled by a tuning parameter  $\epsilon$ . The advantage of our method is that smaller values of  $\epsilon$  can be achieved for the same computational cost, resulting in more accurate inference. We have shown this is the case through asymptotics (Section 3.3) and empirically (Sections 4 and 5.) This increased accuracy allows higher dimensional data or summary statistics to be analysed in practice.

### 6.1 Latent variable considerations

Our method represents the model of interest with latent variables  $x$ , and uses SMC and slice sampling to search for promising  $x$  values. For this search strategy to work well it seems necessary that:

- Evaluating  $y(\theta, x)$  is reasonably cheap.
- Sets of the form  $\{x | d(y_{\text{obs}}, y(\theta, x)) \leq \epsilon\}$  are easy to explore using slice sampling. This would be difficult for sets made up of many disconnected components, or which are lower dimensional manifolds. Smoothness of  $y$  to changes in  $x$  will help meet this condition.

Furthermore, our current implementation requires that the number of latent variables is fixed. However the method could be adapted to more general situations by altering the slice sampling algorithm (see Section 4.2 of Murray and Graham, 2016).

### 6.2 Adaptive and non-adaptive algorithms

The RE-ABC algorithm can use RE-SMC with a fixed  $\epsilon$  sequence (Algorithm 2) or one that is chosen adaptively (Algorithm 3). The non-adaptive RE-SMC algorithm provides

unbiased estimates of the ABC likelihood, as required by the PMMH algorithm, while the adaptive version has a small bias. In practice we observe no very little difference in the posterior results between the two algorithms, suggesting that this bias has a negligible effect in practice. We also note that, if desired, a bias correction approach from C  rou et al. (2012) could be applied.

Nonetheless, we recommend using the non-adaptive RE-SMC algorithm within RE-ABC (together with a pilot run of adaptive RE-SMC to choose the  $\epsilon$  sequence.) The main reason is that it is faster to run in practice, as found in Section 5. Figure 5 shows that this is because the adaptive RE-SMC can terminate more quickly for poor proposed  $\theta$  values. Interestingly, when early termination is not required the adaptive RE-SMC algorithm is slightly quicker. We speculate that this is because it often finds a shorter  $\epsilon$  sequence. Furthermore, the theory of C  rou et al. (2012) suggests that adaptive RE-SMC produces less variable ABC likelihood estimates, which would improve PMMH efficiency. Therefore there may be some scope for a more efficient RE-SMC algorithm which combines the best features of the adaptive and non-adaptive approaches.

### 6.3 Possible extensions

**Joint exploration of  $(\theta, x)$**  RE-ABC proposes many  $\theta$  values which are rejected after calculating an expensive likelihood estimate. An appealing alternative is to update the parameters  $\theta$  conditional on sampled  $x$  values, for example through a Gibbs sampler with state  $(\theta, x)$ . Unfortunately in exploratory analyses of such methods we found the  $\theta$  updates generally did not mix well. The reason is that  $x$  is much more informative for  $\theta$  than the observations  $y_{\text{obs}}$ . This results in small  $\theta$  moves compared to the posterior’s scale. More sophisticated approaches to learning  $\theta$  and  $x$  jointly would be very useful.

**Discrete data** Both RE-SMC algorithms can struggle if there is a discrete data variable  $x^*$ . It can be hard for SMC to move from accepting a set of latent variables  $A$  to another  $A'$

in which the range of possible  $x^*$  values is smaller, because  $\Pr(x \in A' | x \in A, \theta)$  may be very small. The issue is particularly obvious for adaptive RE-SMC as the  $\epsilon$  sequence may fail to fall below some threshold for a large number of iterations. For non-adaptive RE-SMC it would instead result in high-variance likelihood estimates. In Section 5.2 this problem occurs for  $\nu$ , the number of removals. There we adopt an application-specific solution by introducing continuous latent variables (pressure thresholds) into the distance function (5). It would be useful to investigate more general solutions from the rare event literature (e.g. Walter, 2015). We note that despite these potential issues, RE-ABC sometimes performs well with discrete data, as in the binned data model of Section 5.3.

**Non-uniform ABC kernels** In this paper, the ABC likelihood (2) is a convolution of the exact likelihood and a uniform kernel (4). Alternative kernel functions have also been used in ABC (e.g. Wilkinson, 2013) such as a Gaussian:  $k(y; \epsilon) \propto \exp[-\frac{1}{2}d(y, y_{\text{obs}})]$ . The RE-ABC algorithm could easily be adapted to make use of these, but it is not clear what effect it would have on our asymptotic results.

**Estimating log-likelihood gradients** Where log-likelihood gradients can be estimated they allow more efficient inference schemes based on stochastic gradient descent (Poyiadjis et al., 2011) or MCMC (Dahlin et al., 2015). Estimating such gradients from SMC algorithms is possible using the Fisher identity (Poyiadjis et al., 2011). However the calculation will necessarily involve evaluating  $\frac{\partial}{\partial \theta} y(\theta, x)$ , which may be demanding for complicated  $y$  functions. Moreno et al. (2016) use automatic differentiation to evaluate this for some models. Alternatively, Andrieu et al. (2012) propose using infinitesimal perturbation analysis methods. It would be interesting to use either approach with RE-ABC.

**Model choice** A desirable extension to RE-ABC would be methods of model choice. Possible methods to extend our PMMH approach include reversible jump MCMC or using a deviance information criterion. See Chkrebtii et al. (2015) and François and Laval (2011)

for versions of these methods in the ABC context. Alternatively, it may be more fruitful to use our likelihood estimate in algorithms which directly output model evidence estimates, such as importance sampling or population Monte Carlo (Cappé et al., 2004).

**Acknowledgements** We thank Chris Sherlock for suggesting the use of slice sampling and Andrew Golightly for helpful discussions.

## A Computational cost

This appendix justifies the computational costs of ABC and RE-ABC stated in Section 3.3. The argument for ABC is rigorous, while that for RE-ABC is more heuristic. Note that throughout this appendix there is no need to distinguish between the adaptive and non-adaptive versions of RE-ABC.

The results are for the asymptotic regime of small  $\epsilon$  and hold for almost all  $y_{\text{obs}}$ . We make several assumptions:

- A1 The density  $\pi(y|\theta)$  is with respect to Lebesgue measure  $dy$  of dimension  $D$ .
- A2 The distance function is Euclidean distance.
- A3 Running slice sampling once requires  $O(1)$  function evaluations.
- A4 RE-SMC uses  $O(-\log \Pr(d(y, y_{\text{obs}}) \leq \epsilon|\theta))$  iterations.
- A5 The time required to evaluate  $y(\theta, x)$  is bounded above and below by non-zero constants which do not depend on  $\theta$  or  $x$ .

Also, we will usually focus on the case where  $D$  is asymptotically large.

Informally, A1 requires that all components of  $y$  have continuous distributions. Under A2 a key mathematical result below, (6), follows easily. Also, a consequence of A2 which we will use is that, from (3),  $V(\epsilon) \propto \epsilon^{-D}$ . A3 states that the cost of slice sampling does



not increase as  $\epsilon$  shrinks. This is plausible due to our adaptive choice of  $w$  (see Section 2.3), and was empirically verified above (see Figure 3.) It follows that running RE-SMC requires  $O(NT)$  function evaluations: the number is asymptotic to the number of particles multiplied by the number of SMC iterations. A4 states that the number of iterations used by RE-SMC is asymptotically proportional to the log of the rare probability being estimated. This follows from a result of Cérou et al. (2012), reviewed in Section 2.2, that when the RE-SMC algorithm is tuned optimally  $\Pr(A_{k+1}|\theta, x \in A_k)$  is constant, say  $\alpha$ , where  $A_k$  denotes the event  $d(y(\theta, x), y_{\text{obs}}) \leq \epsilon_k$ . Therefore  $\Pr(d(y, y_{\text{obs}}) \leq \epsilon|\theta) = \alpha^T$ , and taking logs gives A4. So the assumption is that RE-SMC is tuned to perform similarly to optimal tuning. Assumption A5 states that performing a simulation has a minimum and maximum time requirement regardless of the inputs, which is usually reasonable. This ensures that computation time is asymptotic to the number of simulations performed.

Many of these assumptions can be weakened. This is discussed in supplementary material, especially for the case of the epidemic model of Section 5.

## A.1 ABC

Consider the probability of a simulation being accepted given  $\theta$ :

$$\Pr(d(y, y_{\text{obs}}) \leq \epsilon|\theta) = \int \pi(y|\theta) \mathbb{1}(d(y, y_{\text{obs}}) \leq \epsilon) dy.$$

By the Lebesgue differentiation theorem (see Stein and Shakarchi, 2009 for example) for almost all  $y_{\text{obs}}$ :

$$\lim_{\epsilon \rightarrow 0} V(\epsilon)^{-1} \int \pi(y|\theta) \mathbb{1}(d(y, y_{\text{obs}}) \leq \epsilon) dy = \pi(y_{\text{obs}}|\theta), \quad (6)$$

Hence for small  $\epsilon$ :

$$\Pr(d(y, y_{\text{obs}}) \leq \epsilon|\theta) \sim V(\epsilon), \quad (7)$$

where  $\sim$  represents an asymptotic relation. (Note that while  $\pi(y_{\text{obs}}|\theta)$  does not affect this asymptotic relationship, the acceptance probability will decrease for small  $\pi(y_{\text{obs}}|\theta)$  i.e. for poor  $\theta$  choices.)

By assumption A5 the time per accepted sample is asymptotic to the number of simulations per accepted sample. Using (7), the latter is asymptotic to  $1/V(\epsilon)$ . In the case of large  $D$  assumption A2 gives that this is  $O(\tau^D)$ , where  $\tau = 1/\epsilon$ . For ABC versions of MCMC and SMC, time per accepted sample (or effective sample) is also bounded below by  $\min_{\theta} \Pr(d(y, y_{\text{obs}}) \leq \epsilon|\theta)^{-1}$ , so the same result applies.

## A.2 RE-ABC

For simplicity we analyse RE-ABC without the possibility of early termination in the RE-SMC algorithm. An algorithm including early termination will give the same output for a smaller computational cost, although we suspect the gain is only likely to be a  $O(1)$  factor. Using the asymptotic results reviewed in Section 2 on SMC likelihood estimation and PMMH we conclude the following. The number of particles in RE-SMC should be  $N = O(T)$  to give a likelihood estimator whose log has variance  $O(1)$ , which optimises efficiency when these estimates are used in PMMH. So, using A3, the number of simulations required by an iteration of RE-ABC is  $O(T^2)$ . Using A4 and (7) gives  $T = O(-\log V(\epsilon))$ .

So the number of simulations required per iteration of RE-ABC is  $O([\log V(\epsilon)]^2)$ . In the case of large  $D$  using A2 gives that this is  $O(D^2[\log \tau]^2)$ . As in the previous section, assumption A5 implies these expressions also give the time per sample of RE-ABC. They are also valid for the more relevant quantity of time per effective sample since effective sample size is proportional to the actual sample size.

## References

- Alquier, P., Friel, N., Everitt, R., and Boland, A. (2016). Noisy Monte Carlo: Convergence of Markov chains with approximate transition kernels. *Statistics and Computing*, 26(1):29–47.
- Andersson, H. and Britton, T. (2000). *Stochastic Epidemic Models and Their Statistical Analysis*. Springer-Verlag.
- Andrieu, C., Doucet, A., and Lee, A. (2012). Contribution to the discussion of Fearnhead and Prangle (2012). *Journal of the Royal Statistical Society: Series B*, 74:451–452.
- Andrieu, C. and Roberts, G. O. (2009). The pseudo-marginal approach for efficient Monte Carlo computations. *The Annals of Statistics*, pages 697–725.
- Barber, S., Voss, J., and Webster, M. (2015). The rate of convergence for approximate Bayesian computation. *Electronic Journal of Statistics*, 9:80–105.
- Beaumont, M. A., Zhang, W., and Balding, D. J. (2002). Approximate Bayesian computation in population genetics. *Genetics*, 162:2025–2035.
- Biau, G., Cérou, F., and Guyader, A. (2015). New insights into approximate Bayesian computation. *Annales de l’Institut Henri Poincaré (B) Probabilités et Statistiques*, 51(1):376–403.
- Blum, M. G. B. (2010). Approximate Bayesian computation: A nonparametric perspective. *Journal of the American Statistical Association*, 105(491):1178–1187.
- Blum, M. G. B., Nunes, M. A., Prangle, D., and Sisson, S. A. (2013). A comparative review of dimension reduction methods in approximate Bayesian computation. *Statistical Science*, 28:189–208.
- Cappé, O., Guillin, A., Marin, J.-M., and Robert, C. P. (2004). Population Monte Carlo. *Journal of Computational and Graphical Statistics*, 13(4):907–929.

- C  rou, F., Del Moral, P., Furon, T., and Guyader, A. (2012). Sequential Monte Carlo for rare event estimation. *Statistics and Computing*, 22(3):795–808.
- Chkrebtii, O. A., Cameron, E. K., Campbell, D. A., and Bayne, E. M. (2015). Trans-dimensional approximate Bayesian computation for inference on invasive species models with latent variables of unknown dimension. *Computational Statistics & Data Analysis*, 86:97–110.
- Dahlin, J., Lindsten, F., and Sch  n, T. B. (2015). Particle Metropolis–Hastings using gradient and Hessian information. *Statistics and computing*, 25(1):81–92.
- Doucet, A., Pitt, M. K., Deligiannidis, G., and Kohn, R. (2015). Efficient implementation of Markov chain Monte Carlo when using an unbiased likelihood estimator. *Biometrika* (online preview).
- Fearnhead, P. and Prangle, D. (2012). Constructing summary statistics for approximate Bayesian computation: Semi-automatic ABC. *Journal of the Royal Statistical Society, Series B*, 74:419–474.
- Forneron, J.-J. and Ng, S. (2015). A likelihood-free reverse sampler of the posterior distribution. *arXiv preprint arXiv:1506.04017*.
- Fran  ois, O. and Laval, G. (2011). Deviance information criteria for model selection in approximate Bayesian computation. *Statistical Applications in Genetics and Molecular Biology*, 10(1).
- Geyer, C. J. (1992). Practical Markov chain Monte Carlo. *Statistical Science*, pages 473–483.
- Graham, M. M. and Storkey, A. (2016). Asymptotically exact conditional inference in deep generative models and differentiable simulators. *arXiv preprint arXiv:1605.07826*.
- Jasra, A. (2015). Approximate Bayesian computation for a class of time series models. *International Statistical Review*.

- Kypraios, T., Neal, P., and Prangle, D. (2016). A tutorial introduction to Bayesian inference for stochastic epidemic models using Approximate Bayesian Computation. *Mathematical Biosciences (online preview)*.
- L’Ecuyer, P., Demers, V., and Tuffin, B. (2007). Rare events, splitting, and quasi-Monte Carlo. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 17(2):9.
- Marin, J.-M., Pudlo, P., Robert, C. P., and Ryder, R. J. (2012). Approximate Bayesian computational methods. *Statistics and Computing*, 22(6):1167–1180.
- Marjoram, P., Molitor, J., Plagnol, V., and Tavaré, S. (2003). Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328.
- McKinley, T. J., Ross, J. V., Deardon, R., and Cook, A. R. (2014). Simulation-based Bayesian inference for epidemic models. *Computational Statistics and Data Analysis*, 71:434–447.
- Meeds, T. and Welling, M. (2015). Optimization Monte Carlo: Efficient and embarrassingly parallel likelihood-free inference. In *Advances in Neural Information Processing Systems*, pages 2071–2079.
- Moreno, A., Adel, T., Meeds, E., Rehg, J. M., and Welling, M. (2016). Automatic variational ABC. *arXiv preprint arXiv:1606.08549*.
- Murray, I. and Graham, M. M. (2016). Pseudo-marginal slice sampling. *Journal of Machine Learning Research*, 51:911–919.
- Neal, P. (2012). Efficient likelihood-free Bayesian computation for household epidemics. *Statistics and Computing*, 22(6):1239–1256.
- Neal, P. and Roberts, G. (2005). A case study in non-centering for data augmentation: stochastic epidemics. *Statistics and Computing*, 15(4):315–327.

- Nott, D. J., Fan, Y., Marshall, L., and Sisson, S. A. (2014). Approximate Bayesian computation and Bayes linear analysis: Toward high-dimensional ABC. *Journal of Computational and Graphical Statistics*, 23(1):65–86.
- Pitt, M. K., Silva, R. D. S., Giordani, P., and Kohn, R. (2012). On some properties of Markov chain Monte Carlo simulation methods based on the particle filter. *Journal of Econometrics*, 171(2):134–151.
- Poyiadjis, G., Doucet, A., and Singh, S. S. (2011). Particle approximations of the score and observed information matrix in state space models with application to parameter estimation. *Biometrika*, 98(1):65–80.
- Prangle, D. (2015). Summary statistics in approximate Bayesian computation. *arXiv preprint arXiv:1512.05633*.
- Sellke, T. (1983). On the asymptotic distribution of the size of a stochastic epidemic. *Journal of Applied Probability*, 20:390–394.
- Sherlock, C., Thiery, A. H., Roberts, G. O., and Rosenthal, J. S. (2015). On the efficiency of pseudo-marginal random walk Metropolis algorithms. *The Annals of Statistics*, 43(1):238–275.
- Sisson, S. A., Fan, Y., and Tanaka, M. M. (2009). Correction: Sequential Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 106(39):16889–16890.
- Smith, R. L. (1996). The hit-and-run sampler: a globally reaching Markov chain sampler for generating arbitrary multivariate distributions. In *Proceedings of the 28th conference on Winter simulation*, pages 260–264. IEEE Computer Society.
- Stein, E. M. and Shakarchi, R. (2009). *Real analysis: measure theory, integration, and Hilbert spaces*. Princeton University Press.

- Streftaris, G. and Gibson, G. J. (2012). Non-exponential tolerance to infection in epidemic systemsmodeling, inference, and assessment. *Biostatistics*, 13(4):580–593.
- Walter, C. (2015). Rare event simulation and splitting for discontinuous random variables. *ESAIM: Probability and Statistics*, 19:794–811.
- Wilkinson, R. D. (2013). Approximate Bayesian computation (ABC) gives exact results under the assumption of model error. *Statistical applications in genetics and molecular biology*, 12(2):129–141.

# A rare event approach to high dimensional Approximate Bayesian computation Supplementary material

Dennis Prangle, Richard G. Everitt and Theodore Kypraios

## 1 Further results from Abakaliki example

This section reports some further details of our analysis of the Abakaliki data. Table 1 contains parameter estimates from adaptive RE-ABC analyses. Figure 1 shows simulated epidemics from each model using RE-ABC. Figure 2 shows trace plots of log-likelihood estimates produced by RE-ABC. These are all discussed in the main text.

Model	$R_0$	Pressure thresholds		Infectious period	
		Mean	Standard deviation	Mean	Standard deviation
5 day bins	1.17 (0.29)	0.11 (0.03)	0.11 (0.03)	11.4 (2.9)	11.4 (2.9)
Gamma infectious period	1.16 (0.22)	0.08 (0.03)	0.08 (0.03)	14.9 (4.2)	6.2 (1.8)
Weibull pressure thresholds	-	0.10 (0.04)	0.11 (0.04)	12.2 (3.8)	12.2 (3.8)

Table 1: Adaptive RE-ABC posterior estimates of basic reproduction number  $R_0$  and the means and standard deviations of pressure thresholds and infectious periods for the Abakaliki data under three models. Monte Carlo estimates are quoted along with standard deviations in brackets. The  $R_0$  value is not given for the Weibull pressure threshold model as no definition is available for this model.



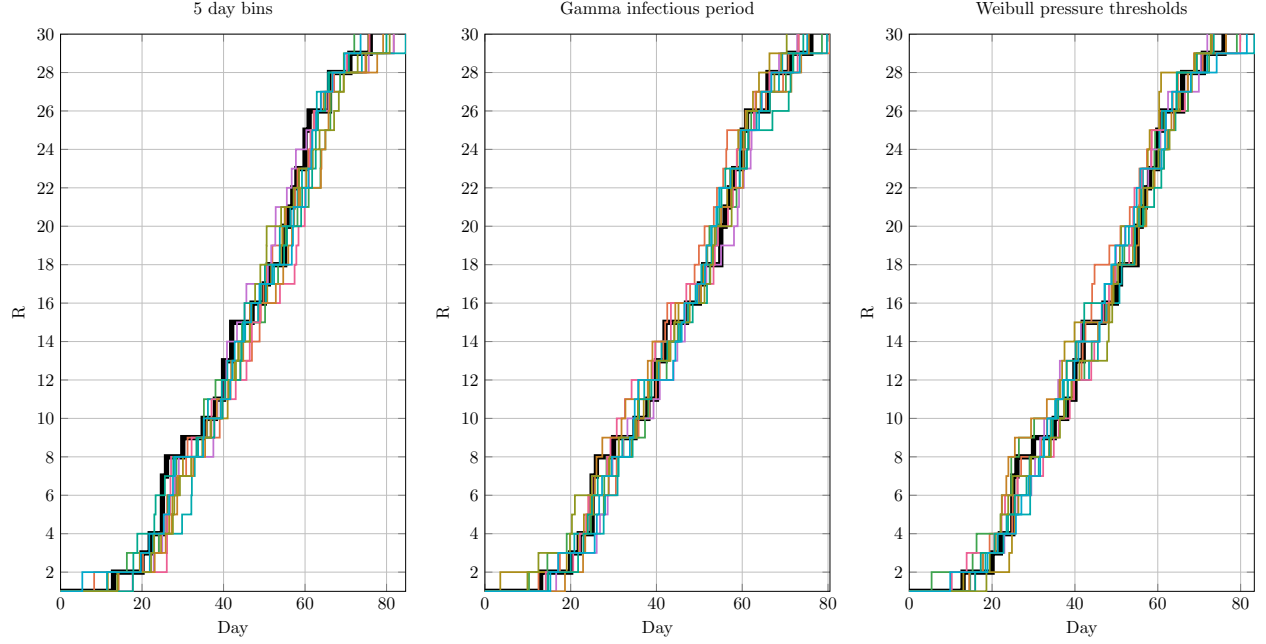


Figure 1: Observed data (black line) and accepted simulations under RE-ABC (coloured lines) for an analysis of Abakaliki data. The x-axis shows number of days since the first removal, and the y-axis shows the total number of removals so far. Accepted simulations are generated by running RE-ABC for 10 parameter values taken from thinned MCMC output and selecting one particle from the final iteration whose distance is below the threshold  $\epsilon$ .

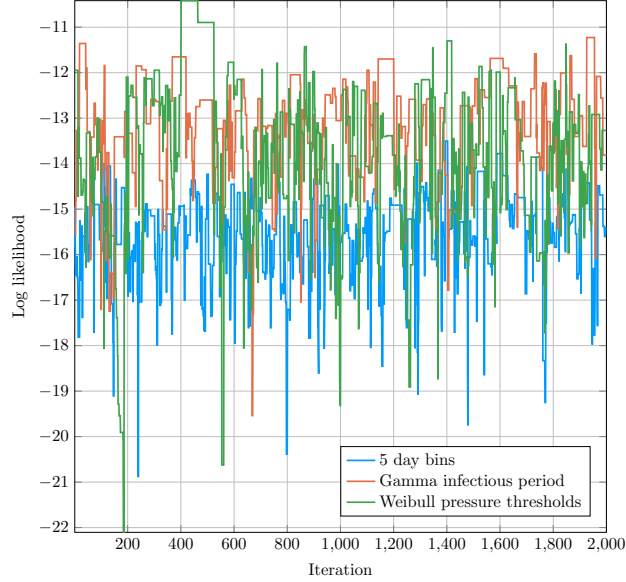


Figure 2: Trace plots of log likelihood estimates resulting from RE-ABC analyses of Abakaliki data.

## 2 Asymptotics with weaker assumptions

This section discusses extending the asymptotic theory of the main paper to use weaker assumptions, in particular showing how it can be used with most of the SIR models from the main paper. For this model assumption A1 (data has density with respect to Lebesgue measure) does not hold because the data is not continuous, instead involving a discrete observation of the number of removals,  $\nu$ , and  $\nu - 1$  continuous inter-removal times. Furthermore assumption A2 (distance function is Euclidean) does not hold because a more complicated distance function was needed. Our argument can be adapted to this model by showing that both assumptions effectively hold for sufficiently small  $\epsilon$  values. This is discussed in Section 2.1.

A further problem arises because the observed data contains repeated recovery times. This means the data is on the boundary of the model's support, which causes technical problems. While the asymptotic results of the main paper remain true for almost all  $y_{\text{obs}}$  values, they are not necessarily valid when  $y_{\text{obs}}$  is on this boundary. This problem is discussed in Section 2.2.

Finally, note that the main paper includes a model with discrete summary statistics: days of removal rounded down to a multiple of 5. Here a sufficiently small non-zero  $\epsilon$  value ensures an exact match of simulated and observed data. Therefore it is not of interest to consider small  $\epsilon$  asymptotics for this case.

### 2.1 Weakening assumptions A1 and A2

Suppose that assumptions A1 and A2 do not hold, but there is some  $\epsilon_0 > 0$  with the following properties.

B1 There is an injective mapping  $z$  from  $A = \{y | d(y, y_{\text{obs}}) < \epsilon_0\}$  to  $\mathbb{R}^{D'}$ .

B2 The distribution  $z(y) | \theta, y \in A$  has density  $\pi(z | \theta)$  with respect to Lebesgue measure  $dy$  of dimension  $D'$ .

B3 For  $y \in A$ ,  $d(y, y_{\text{obs}})$  equals  $d_E(z(y), z_{\text{obs}})$  where  $d_E$  denotes Euclidean distance and  $z_{\text{obs}} = z(y_{\text{obs}})$ .

### 2.1.1 SIR model

For the SIR model we can select  $\epsilon_0$  such that  $d(y, y_{\text{obs}}) < \epsilon_0$  guarantees that  $y$  has the same number of removals as  $y_{\text{obs}}$ . For example  $\epsilon_0 = k$  will achieve this (recall that  $k$  is a penalty in the distance function for the wrong number of removals). Suppose  $y$  has removal times  $r_{(1)} \leq r_{(2)} \leq \dots \leq r_{(\nu)}$  where  $\nu$  is the number of removal times. Let  $z(y)$  be the dimension  $\nu - 1$  vector of times since first removal i.e. the values  $s_{(i)} = r_{(i)} - r_{(1)}$  for  $2 \leq i \leq \nu$ . This meets the assumptions B1-B3.

### 2.1.2 Asymptotics

For both ABC and RE-ABC the crucial quantity is  $\Pr(d(y, y_{\text{obs}}) \leq \epsilon | \theta)$ . For  $\epsilon < \epsilon_0$  this is given by:

$$\Pr(y \in A | \theta) \Pr(d(y, y_{\text{obs}}) \leq \epsilon | \theta, y \in A)$$

Let these probabilities be  $P_1$  and  $P_2$  respectively. Since the former does not depend on  $\epsilon$  we have  $P_1 = O(1)$ . The latter is

$$\Pr(d_E(z(y), z_{\text{obs}}) \leq \epsilon | \theta, y \in A) = \int \pi(z | \theta) \mathbb{1}(d_E(z, z_{\text{obs}}) \leq \epsilon) dz$$

Now we can repeat the argument of the main paper. By the Lebesgue differentiation theorem for almost all  $z_{\text{obs}}$ :

$$\lim_{\epsilon \rightarrow 0} \frac{\int \pi(z | \theta) \mathbb{1}(d(z, z_{\text{obs}}) \leq \epsilon) dz}{\int \mathbb{1}(d(z, z_{\text{obs}}) \leq \epsilon) dz} = \pi(z_{\text{obs}} | \theta).$$

Hence for small  $\epsilon$ ,

$$P_2 \sim \int \mathbb{1}(d(z, z_{\text{obs}}) \leq \epsilon) dz = O(\epsilon^{D'}).$$

Using the arguments in the main paper it follows that the time per sample in ABC is  $(P_1 P_2)^{-1}$  which is  $O(\epsilon^{-D'})$ , and the time per effective sample for RE-ABC is  $O(D'^2 [\log \epsilon]^2)$ .

## 2.2 Data in the boundary of the support

Our asymptotics rely on the Lebesgue differentiation theorem which states that when  $g(y)$  is Lebesgue integrable and  $dy$  is Lebesgue measure then the following holds for almost all  $y_0$ :

$$\lim_{\epsilon \rightarrow 0} \frac{\int g(y) \mathbb{1}(y \in B_\epsilon) dy}{\int \mathbb{1}(y \in B_\epsilon) dy} = g(y_0), \quad (1)$$

where  $B_\epsilon$  is a ball of radius  $\epsilon$  centred on  $y_0$ .

However this is not true for  $y_0$  on the boundary of the support of  $g(y)$ . For example suppose  $g(y)$  is a uniform density on  $[0, 1]$  and  $y_0 = 0$ . Then for  $\epsilon < 1$ :

$$\frac{\int g(y) \mathbb{1}(y \in B_\epsilon) dy}{\int \mathbb{1}(y \in B_\epsilon) dy} = 1/2$$

This problem can be avoided by replacing  $B_\epsilon$  with  $B'_\epsilon = B_\epsilon \cap \text{supp}(g)$ . The Lebesgue differentiation theorem remains true in this case (see Stein and Shakarchi, 2009), as the  $B'_\epsilon$  sets meet the condition of *bounded eccentricity*. That is, each  $B'_\epsilon$  is contained in some ball  $B$  such that  $|B'_\epsilon| \leq \delta |B|$  for some constant  $\delta$ . (This method could be also be used to show our asymptotics hold for many non-Euclidean distance functions.)

## References

Stein, E. M. and Shakarchi, R. (2009). *Real analysis: measure theory, integration, and Hilbert spaces*. Princeton University Press.